# Bounded regret in stochastic multi-armed bandits

SÉBASTIEN BUBECK, VIANNEY PERCHET[*] AND PHILIPPE RIGOLLET[†]

*Princeton University, Université Paris Diderot and Princeton University*

*Abstract.* We study the stochastic multi-armed bandit problem when one knows the value $\mu^{(\star)}$ of an optimal arm, as a well as a positive lower bound on the smallest positive gap $\Delta$. We propose a new randomized policy that attains a regret *uniformly bounded over time* in this setting. We also prove several lower bounds, which show in particular that bounded regret is not possible if one only knows $\Delta$, and bounded regret of order $1/\Delta$ is not possible if one only knows $\mu^{(\star)}$.

## 1. INTRODUCTION

In this paper we investigate the classical stochastic multi-armed bandit problem introduced by [12] and described as follows: an agent facing $K$ actions (or bandit arms) selects one arm at every time step until a finite time horizon $n \geq 1$. Successive pulls of each arm $i \in \{1, \ldots, K\}$ yield a sequence of i.i.d rewards $Y_1^{(i)}, Y_2^{(i)}, \ldots$ according to some unknown distribution $\nu_i$ with expected value $\mu^{(i)}$. Denote by $\star \in \{1, \ldots, K\}$ any optimal arm defined such that $\mu^{(\star)} = \max_{i=1,\ldots,K} \mu^{(i)}$. A *policy* $I = \{I_t\}$ is a sequence of random variables $I_t \in \{1, ..., K\}$ indicating which arm to pull at each time $t = 1, \ldots, n$ and such that $I_t$ depends only on observations strictly anterior to $t$. The performance of a policy $I$ is measured by its (cumulative) *regret* at time $n$ that is defined by

$$R_n = n\mu^{(\star)} - \sum_{t=1}^n \mathbb{E}\,\mu^{(I_t)}\,.$$

Observe that if we denote by $T_i(t) = \sum_{\ell=1}^{t-1} \mathbb{1}\{I_\ell = i\}$ the number of times arm $i$ was pulled (strictly) before time $t \geq 2$ and by $\Delta_i = \mu^{(\star)} - \mu^{(i)}$ the gap between arm $i$ and the optimal arm, then one can rewrite the regret as $R_n = \sum_{i=1}^K \Delta_i \mathbb{E}T_i(n+1)$. This formulation will be used hereafter.

We refer the reader to [5] for a survey of the extensive literature on this problem and its variations. In this paper we investigate a phenomenon that was first observed in [8]: with some prior knowledge (in the form of lower bounds) on

---

the maximal mean $\mu^{(\star)}$ and the minimal gap $\Delta = \min_{i:\Delta_i>0} \Delta_i$, it is possible to obtain a regret that is *bounded uniformly in $n$*, which implies in particular that the regret does not tend to infinity as the time horizon $n$ tends to infinity. Note that this result is striking, as the seminal paper [9] indicates that, if one has no prior knowledge on the distributions, then asymptotically (in $n$) a regret of order $\log n$ is unavoidable.

### 1.1 Contributions

We describe in Section 2 a simple algorithm for the two-armed bandit problem when one knows the largest expected reward $\mu^{(\star)}$ and the gap $\Delta$. In this two-armed case, this amounts to knowing $\mu^{(1)}$ and $\mu^{(2)}$ up to a permutation. We show that the regret of this algorithm is bounded by $\Delta + 16/\Delta$, uniformly in $n$. The optimality of this bound is assessed in Section 4 where we show that any agent knowing $\Delta$ and $\mu^{(\star)}$ must incur a regret of at least $1/\Delta$. This upper and lower bounds raise the following question: can such bounded regret be achieved without one of these two pieces of information? It follows from Theorems 6 and 8 that the answer to this question is negative. Indeed, the sole knowledge of either $\Delta$ or $\mu^{(\star)}$ leads to a rescaled regret $\Delta R_n$ that is at least logarithmic in $n$. Interestingly, all these results are fully non-asymptotic, including lower bounds.

What if $\Delta$ is not perfectly known but only $\varepsilon > 0$ such that $\Delta > \varepsilon$? We answer this question in Section 3 in the context of the general $K$-armed bandit problem. There, we prove an upper bound on $R_n$ when one knows the maximal mean $\mu^{(\star)}$ together with a positive lower bound $\varepsilon$ on the smallest gap $\Delta$. Specifically, we design a randomized policy for which

$$R_n \le \sum_{i:\Delta_i>0} \left\{ \Delta_i + \frac{32}{\Delta_i} \log\left(\frac{5}{\varepsilon}\right) \right\}.$$

Moreover, it follows form our main lower bound in Theorem 8 that this result cannot be improved without further assumptions, since for $\varepsilon$ of order of $1/\sqrt{n}$ —no information on the smallest gap— a logarithmic growth in $n$ is unavoidable for the rescaled regret $\Delta R_n$. However for $\varepsilon$ of order $\Delta$ one would expect no dependency on $\varepsilon$ (since at least for $K = 2$ our policy of Section 2 attains a regret of order $1/\Delta$). To deal with this issue we propose an improvement of the basic policy that for which the term $\log(1/\varepsilon)$ is replaced by $\log(\Delta_i/\varepsilon) \log \log \varepsilon$. In particular if all the gaps $\Delta_i$ and $\varepsilon$ are of the same order, the logarithmic becomes a log-log term.

The *exploration-exploitation tradeoff* is a preponderant paradigm in the bandit literature. The effects of this tradeoff already appear for the case $K = 2$ in the form of the $\log n$ term derived in the original [9] paper. Indeed, there exist simple classes of (two!) problems over which the regret is uniformly bounded with full information but cannot be bounded uniformly with bandit feedback, see Theorem 6. Clearly, this tradeoff should become more and more apparent as the number of arms increases but this is not our main focus. Rather, the combination of our results sheds light on an interesting phenomenon: the effects of the tradeoff vanish when both $\Delta$ and $\mu^{(\star)}$ are known but can be seen already when $K = 2$ and either $\Delta$ or $\mu^{(\star)}$ is unknown.

### 1.2 Related works

The two-armed bandit problem when one knows the distributions of the arms up to a permutation was first investigated in [8]. The authors observed that in that case, using a policy based on the sequential likelihood ratio test, one can obtain a regret uniformly bounded over $n$. Both upper and lower bounds were provided. This setting was generalized in [7], where the authors considered the general multi-armed bandit problem when one knows a separating value $\gamma$ between the largest mean and the other means. In that case they proved the bounded regret property for a policy based on sequential likelihood ratio tests for $H_0 : \mu > \gamma$ vs. $H_1 : \mu < \gamma$ (assuming exponential distributions to compute the likelihoods). They also designed a more subtle strategy for the case when only $\mu^{(\star)}$ is known. In that case too they proved a bounded regret property. The main open problems left by these works are (i) to understand the limitations of bounded regret, and (ii) to characterize the exact dependence on the parameters in the regret (when bounded regret is achievable). In this paper we make progress on both questions.

Regarding the limitations of bounded regret, we prove three finite-time lower bounds, including a finite-time version of the seminal result of [9]. Ideas similar to the ones we develop in Theorems 5 and 6 already appeared in [6] but our results are fully non asymptotic with the exact dependence in the parameters involved. Theorem 8 is more innovative. It shows that a logarithmic growth for the rescaled regret $\Delta R_n$ is unavoidable even if one knows $\mu^{(\star)}$. The proof of this result goes beyond any previous lower bound for the stochastic multi-armed bandit problem, including [7, 9], since all of them required to distinguish problems with different values of $\mu^{(\star)}$ (such as the ones in Theorem 6 for example). As a consequence of this theorem, we can deduce that the policies with bounded regret derived in [7, 1] with only the knowledge of $\mu^{(\star)}$ must have a suboptimal dependency in $1/\Delta$.

The knowledge of $\mu^{(\star)}$ was also exploited in other works. For instance in [13], the authors showed that knowing $\mu^{(\star)}$ allows for policies with provably better concentration properties. Their policies are based on sequential likelihood ratio tests for $H_0 : \mu = \mu^{(\star)}$ vs. $H_1 : \mu < \mu^{(\star)}$ (assuming Gaussian distributions to compute the likelihoods). To some extent it was to be expected that the knowledge of $\mu^{(\star)}$ leads to an improved regret as it partially removes the need for exploration: if one arm has empirical performances close to $\mu^{(\star)}$, one can be confident that this is the best arm without worrying that it could be the best arm only because we have not yet explored enough the other options. However note that the problem turns out to be more subtle than the above simple argument and underlines the fact that one needs more than the knowledge of $\mu^{(\star)}$ in order to have a bounded regret with optimal scaling in $1/\Delta$. Indeed, Theorem 8 implies that the sole knowledge of $\mu^{(\star)}$ does not warrant the bounded property for the rescaled regret $\Delta R_n$.

### 1.3 Basic assumptions

Throughout the paper, we assume that the distributions $\nu_i$ are sub-Gaussian that is $\int e^{\lambda(x-\mu)}\nu_i(dx) \leq e^{\lambda^2/2}$ for all $\lambda \in \mathbb{R}$. Note that these include Gaussian distributions with variance less than 1 and distributions supported on an interval of length less than 2.

We denote by $\widehat{\mu}_s^{(i)} = \frac{1}{s}\sum_{\ell=1}^s Y_\ell^{(i)}$ the empirical mean of arm $i$ after $s$ pulls,

for $s \geq 1$. Together with a Chernoff bound, it is not hard to see that the sub-Gaussian assumption implies the following concentration inequality, valid for any $u > 0$,

$$(1.1) \qquad \mathbb{P}(\widehat{\mu}_s^{(i)} - \mu^{(i)} > u) \leq \exp\left(-\frac{su^2}{2}\right).$$

## 2. THE TWO-ARMED CASE

In this section we investigate a toy example where $K = 2$ and the agent knows exactly both $\mu^{(\star)} = 0$ (without loss of generality) and $\Delta$. While somewhat simplistic this example offers a convenient framework to lay the main ideas to build policies with bounded regret.
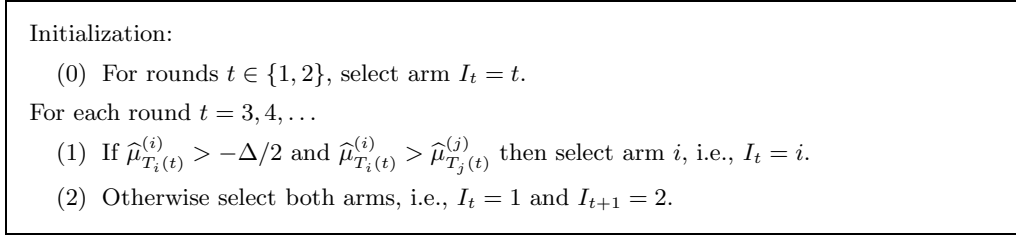
---

Initialization:

(0) For rounds $t \in \{1, 2\}$, select arm $I_t = t$.

For each round $t = 3, 4, \ldots$

(1) If $\widehat{\mu}_{T_i(t)}^{(i)} > -\Delta/2$ and $\widehat{\mu}_{T_i(t)}^{(i)} > \widehat{\mu}_{T_j(t)}^{(j)}$ then select arm $i$, i.e., $I_t = i$.

(2) Otherwise select both arms, i.e., $I_t = 1$ and $I_{t+1} = 2$.

---

FIGURE 1. *A policy with bounded regret for the two-armed bandit problem.*

THEOREM 1.  *Policy 1 has regret bounded as $R_n \leq \Delta + 16/\Delta$, uniformly in* $n$.

PROOF. Without loss of generality we assume that $1 = \star$ is the optimal arm. Observe that

$$\{I_t = 2\} \subset \{t = 2\} \cup \{\widehat{\mu}_{T_2(t)}^{(2)} > -\Delta/2, t \geq 3, \ I_t = 2\} \cup \{\widehat{\mu}_{T_2(t)}^{(2)} \leq -\Delta/2, t \geq 3, \ I_t = 2\}.$$

Summing over $t$ for the second event, we get
$$(2.2)$$
$$\mathbb{E} \sum_{t=3}^{n} \mathbb{1}\{\widehat{\mu}_{T_2(t)}^{(2)} > -\Delta/2, \ I_t = 2\} \leq \mathbb{E} \sum_{t=1}^{n} \mathbb{1}\{\widehat{\mu}_t^{(2)} > -\Delta/2\} \leq \sum_{t=1}^{n} \exp(-t\Delta^2/8) \leq \frac{8}{\Delta^2}.$$

For the third event we use the definition of the policy to obtain

$$\{\widehat{\mu}_{T_2(t)}^{(2)} \leq -\Delta/2, t \geq 3, \ I_t = 2\} \subset \{\widehat{\mu}_{T_1(t-1)}^{(1)} \leq -\Delta/2, t \geq 3, \ I_{t-1} = 1\}$$

and conclude as in (2.2).                                           □

This policy has two weaknesses. First one may pay a big price for misspecifying the value of $\Delta$. Namely if one only knows a lower bound $0 < \varepsilon \leq \Delta$ and substitutes $\varepsilon$ to $\Delta$ in Policy 1, then it follows easily that the regret becomes of order $\Delta/\varepsilon^2$. Furthermore, for essentially the same reason, the trivial generalization of this algorithm to the $K$-armed case would give a regret bounded by $\sum_i \Delta_i/\Delta^2$. In the next section we show how to overcome these two issues using a new, randomized, policy.

## 3. A FAMILY OF POLICIES WITH BOUNDED REGRET

In this section we consider the general multi-armed case, when the agent knows $\mu^{(\star)} = 0$ (without loss of generality) and an $\varepsilon > 0$ such that $\varepsilon \leq \Delta$. Akin to Policy 1, the policy analyzed here sets a threshold at $-\varepsilon/2$ and prescribes to pull a single arm above this threshold. However if all arms have their empirical mean below this threshold, then the policy is more subtle than what was described in the previous section (where all arms were pulled in round robin fashion). Here the policy picks an arm at random, where the probability of selecting arm $i$ is essentially proportional to $(\hat{\mu}_{T_i(t)}^{(i)})^{-2}$, which is an empirical estimate of $\Delta_i^{-2}$ since $\mu^{(\star)} = 0$. Policy 2 is slighly more general, as it uses a potential function $\psi : \mathbb{R}_+ \to \mathbb{R}_+$, and selects arm $i$ with probability inversely proportional to $\psi(|\hat{\mu}_{T_i(t)}^{(i)}|)$. The natural choice is $\psi(x) = x^2$, but other choices can lead to improved performances, see Theorem 2 below. Note that we also analyze the case where $\varepsilon = 0$ (that is, when we have no information on the smallest gap).
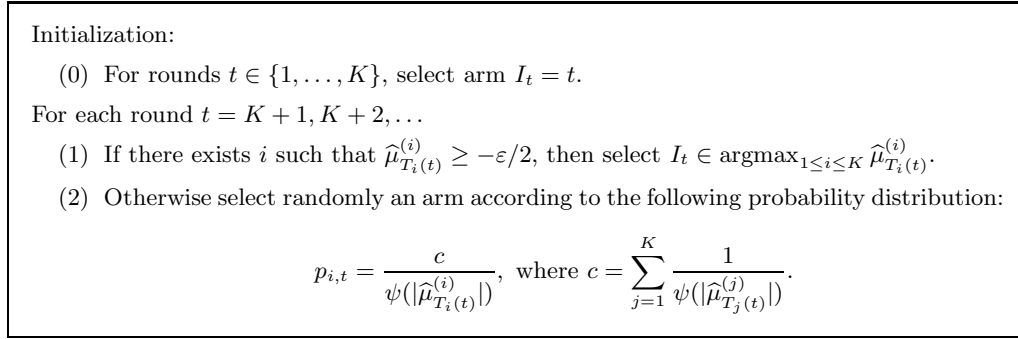
---

Initialization:

(0) For rounds $t \in \{1, \ldots, K\}$, select arm $I_t = t$.

For each round $t = K + 1, K + 2, \ldots$

(1) If there exists $i$ such that $\hat{\mu}_{T_i(t)}^{(i)} \geq -\varepsilon/2$, then select $I_t \in \text{argmax}_{1 \leq i \leq K} \hat{\mu}_{T_i(t)}^{(i)}$.

(2) Otherwise select randomly an arm according to the following probability distribution:

$$p_{i,t} = \frac{c}{\psi(|\hat{\mu}_{T_i(t)}^{(i)}|)}, \quad \text{where } c = \sum_{j=1}^{K} \frac{1}{\psi(|\hat{\mu}_{T_j(t)}^{(j)}|)}.$$

---

FIGURE 2. *A family of policies with bounded regret for the $K$-armed bandit problem.*

THEOREM 2. *Fix $\varepsilon \in (0, 1 \wedge \Delta]$, then Policy 2 associated with the potential $\psi(x) = x^2$ satisfies for all $n \geq 1$,*

$$(3.3) \qquad R_n \leq \sum_{i:\Delta_i > 0} \left\{ \Delta_i + \frac{32}{\Delta_i} \log\left(\frac{5}{\varepsilon}\right) \right\}.$$

*Furthermore for $\varepsilon = 0$, let $v = \mathbb{E}\left(Y_1^{(\star)}\right)^2$, then the regret is bounded as*

$$(3.4) \qquad R_n \leq \sum_{i:\Delta_i > 0} \left\{ \Delta_i + (1 \vee v) \frac{4 \log(9n)}{\Delta_i} \right\}.$$

*The dependency in $\varepsilon$ can be reduced by using the potential $\psi(x) = \frac{x^2}{\log(4x/\varepsilon)}$ since it yields*

$$(3.5) \qquad R_n \leq \sum_{i:\Delta_i > 0} \left\{ \Delta_i + \frac{32 \log\left(\frac{2\Delta_i}{\varepsilon}\right)}{\Delta_i} \left[3 + \log\log\left(\frac{4}{\varepsilon}\right)\right] \right\}.$$

If $\varepsilon$ is of the order of every $\Delta_i$, then Equation (3.5) upper bounds the regret in $\sum_i \log\log(1/\Delta_i)/\Delta_i$; on the other hand, using the potential $\psi(x) = x^2$ only guarantees, under the same assumptions, a bound in $\sum_i \log(1/\Delta_i)/\Delta_i$.

The result for $\varepsilon = 0$ implies that when one has no information on the smallest gap, our policy does not obtain bounded regret but it recovers the performances of UCB, [3]. As we shall see in Section 4 it is in fact impossible to obtain bounded regret scaling in $1/\Delta$ if one only knows $\mu^{(\star)}$.

Theorem 2 is deduced from the following more general regret bound for Policy 2 expressed in terms of the properties of the potential $\psi$.

THEOREM 3.    *Fix* $\varepsilon \in [0, \Delta]$ *and let* $\psi$ *be a differentiable and increasing function* $\psi : [\varepsilon/2, +\infty) \to \mathbb{R}^+$. *If* $\varepsilon > 0$, *Policy 2 satisfies for all* $n \geq 1$,

$$(3.6) \qquad R_n \leq \sum_{i:\Delta_i>0} \left\{ \Delta_i + \frac{8}{\Delta_i} + \frac{\Delta_i}{\psi(\Delta_i/2)} \left[ \frac{8\psi(\varepsilon/2)}{\varepsilon^2} + \int_{\varepsilon/2}^{+\infty} \frac{2\psi'(x)}{e^{\frac{x^2}{2}} - 1} dx \right] \right\}.$$

*Furthermore for* $\varepsilon = 0$ *it satisfies*

$$(3.7) \qquad R_n \leq \sum_{i:\Delta_i>0} \left( \Delta_i + \frac{8}{\Delta_i} + \frac{\Delta_i}{\psi(\Delta_i/2)} \sum_{t=1}^{n} \mathbb{E} \, \psi(|\widehat{\mu}_t^{(1)}|) \right).$$

PROOF. Without loss of generality we assume that $1 = \star$ is the optimal arm. We decompose the event of a wrong selection into three events:

$$\{I_t = i\} \subset \{t = i\} \cup \{\widehat{\mu}_{T_i(t)}^{(i)} > -\Delta_i/2, \, t \geq K + 1, \, I_t = i\}$$
$$\cup \{\widehat{\mu}_{T_i(t)}^{(i)} \leq -\Delta_i/2, \, t \geq K + 1, \, I_t = i\}.$$

Using (2.2) one can easily prove that the cumulative probability of the first two events is smaller than $1 + 8/\Delta_i^2$. For the third event, it is convenient to define the random variable $Z \in \{0, 1, 2\}$ that indicates whether the agent plays according to (0), (1) or (2) in Policy 2. We write the following, using the definition of the algorithm and the fact that $\psi$ is non-decreasing,

$$\mathbb{P}\{\widehat{\mu}_{T_i(t)}^{(i)} \leq -\Delta_i/2, \, t \geq K + 1, \, I_t = i\} = \mathbb{P}\{\widehat{\mu}_{T_i(t)}^{(i)} \leq -\Delta_i/2, \, I_t = i, Z = 2\}$$

$$= \mathbb{E} \, p_{i,t} \mathbb{1}\{\widehat{\mu}_{T_i(t)}^{(i)} \leq -\Delta_i/2, Z = 2\} = \mathbb{E} \, \frac{p_{i,t}}{p_{1,t}} p_{1,t} \mathbb{1}\{\widehat{\mu}_{T_i(t)}^{(i)} \leq -\Delta_i/2, Z = 2\}$$

$$\leq \mathbb{E} \, \frac{\psi(|\widehat{\mu}_{T_1(t)}^{(1)}|)}{\psi(\Delta_i/2)} p_{1,t} \mathbb{1}\{\widehat{\mu}_{T_i(t)}^{(i)} \leq -\Delta_i/2, Z = 2\} \leq \frac{1}{\psi(\Delta_i/2)} \, \mathbb{E} \, \psi(|\widehat{\mu}_{T_1(t)}^{(1)}|) p_{1,t} \mathbb{1}\{Z = 2\}$$

$$\leq \frac{1}{\psi(\Delta_i/2)} \, \mathbb{E} \, \psi(|\widehat{\mu}_{T_1(t)}^{(1)}|) \mathbb{1}\{\widehat{\mu}_{T_1(t)}^{(1)} < -\varepsilon/2, \, t \geq K + 1\}.$$

A simple rewriting of time then concludes the proof for the case of $\varepsilon = 0$. We use the slight abuse of notation $\psi^{-2}(x) := [\psi^1(x)]^2$, and $\psi(\infty) = \lim_{x \to +\infty} \psi(x)$. For $\varepsilon > 0$ we have

$$\sum_{t=1}^{n} \mathbb{E}\, \psi(|\widehat{\mu}_{T_1(t)}^{(1)}|) \mathbb{1}\{\widehat{\mu}_{T_1(t)}^{(1)} \leq -\varepsilon/2\} \leq \sum_{t=1}^{n} \mathbb{E}\, \psi(|\widehat{\mu}_t^{(1)}|) \mathbb{1}\{\widehat{\mu}_t^{(1)} \leq -\varepsilon/2\}$$

$$= \sum_{t=1}^{n} \int_0^{+\infty} \mathbb{P}\left(\psi(|\widehat{\mu}_t^{(1)}|) \mathbb{1}\{\widehat{\mu}_t^{(1)} \leq -\varepsilon/2\} \geq x\right) dx$$

$$= \sum_{t=1}^{n} \left\{ \psi\left(\frac{\varepsilon}{2}\right) \mathbb{P}\left(|\widehat{\mu}_t^{(1)}| > \frac{\varepsilon}{2}\right) + \int_{\psi(\varepsilon/2)}^{\psi(\infty)} \mathbb{P}(\psi(|\widehat{\mu}_t^{(1)}|) \geq x) dx \right\}$$

$$\leq \sum_{t=1}^{n} \left\{ \psi\left(\frac{\varepsilon}{2}\right) e^{-\frac{t\varepsilon^2}{8}} + \int_{\psi(\varepsilon/2)}^{\psi(\infty)} 2e^{-\frac{t\psi^{-2}(x)}{2}} dx \right\}$$

$$\leq \frac{8}{\varepsilon^2} \psi\left(\frac{\varepsilon}{2}\right) + \int_{\psi(\varepsilon/2)}^{\psi(\infty)} \frac{2}{e^{\frac{\psi^{-2}(x)}{2}} - 1} dx.$$

Making the change of variable $x = \psi(u)$ concludes the proof of Theorem 3. $\square$

Theorem 2 follows from Theorem 3 with specific choices for $\psi$. First, take $\psi(x) = x^2$, $\varepsilon \in (0, 1]$ and observe that the integral in (3.6) can be computed as

$$\int_{\varepsilon/2}^{+\infty} \frac{4x}{e^{\frac{x^2}{2}} - 1} dx = -4\log\left(1 - e^{-\frac{\varepsilon^2}{8}}\right) \leq 8\log\left(\frac{3}{\varepsilon}\right),$$

which gives (3.3). When $\varepsilon = 0$, since $\mathbb{E}\, \psi(|\widehat{\mu}_t^{(1)}|) = v/t$, Equation (3.7) directly gives (3.4).

Next, we turn to the the slightly more sophisticated potential function $\psi(x) = \frac{x^2}{\log(4x/\varepsilon)}$. Observe that for any $x \geq 0$,

$$\psi'(x) = \frac{2x}{\log(4x/\varepsilon)} - \frac{x}{\log^2(4x/\varepsilon)} \leq \frac{2x}{\log(4x/\varepsilon)}.$$

Therefore, for $\varepsilon \in (0, 1]$, the integral in (3.6) is bounded from above by

$$\int_{\varepsilon/2}^{+\infty} \frac{4x}{\log(4x/\varepsilon)[e^{\frac{x^2}{2}} - 1]} dx \leq \int_{\varepsilon/2}^{1} \frac{8}{x\log(4x/\varepsilon)} dx + \int_1^{\infty} 9e^{-\frac{x^2}{2}} dx$$

$$\leq 8\log\log(4/\varepsilon) - 8\log\log 2 + 4$$

$$\leq 8\log\log(4/\varepsilon) + 7.$$

It concludes the proof of (3.5).

## 4. LOWER BOUNDS

We conclude our study of bounded regret in stochastic multi-armed bandits with three different lower bounds. For simplicity, we phrase these results for the simple two-armed case. First we show with Theorem 5 that if one knows both $\mu^{(\star)}$ and $\Delta$, then the best attainable regret is of order $1/\Delta$, which matches (up to a numerical constant) the result of Theorem 1. Next we show in Theorem 6 that the sole knowledge of $\Delta$ leads to a lower bound of order $\log(n\Delta^2)/\Delta$. This theorem implies that the bounds of [2], [4] and [10] exhibit a tight dependence in

$\Delta$ (for the two-armed case), unlike the famous result of [9]. Moreover, compared to the proof of [9], our approach is (i) much simpler, (ii) non-asymptotic and (iii) it is not limited to a certain class of policies. Finally we show in Theorem 8 that if one only knows $\mu^{(\star)}$ then a regret of order $\frac{\log(n)}{\Delta}$ is unavoidable (for some value of $\Delta$).

Our proof strategy consists in rephrasing arm selection as a hypothesis testing problem, and then use well-known lower bounding techniques for the minimax risk of hypothesis testing. For instance, the proof of Theorem 5 and Theorem 6 builds upon the following result; see [14, Chaper 2] for a proof, or Lemma 7 below with $\lambda$ chosen to be a Dirac mass at 1. Recall that the Kullback-Leibler divergence between two positive measures $\rho, \rho'$ with $\rho'$ absolutely continuous with respect to $\rho$, is defined as

$$\mathrm{KL}(\rho, \rho') = \int \log\left(\frac{d\rho}{d\rho'}\right) d\rho = \mathbb{E}_{X \sim \rho} \log\left(\frac{d\rho}{d\rho'}(X)\right).$$

LEMMA 4. *Let $\rho_0, \rho_1$ be two probability distributions supported on some set $\mathcal{X}$, with $\rho_1$ absolutely continuous with respect to $\rho_0$. Then for any measurable function $\psi : \mathcal{X} \to \{0, 1\}$, one has*

$$\mathbb{P}_{X \sim \rho_0}(\psi(X) = 1) + \mathbb{P}_{X \sim \rho_1}(\psi(X) = 0) \geq \frac{1}{2}\exp\left(-\mathrm{KL}(\rho_0, \rho_1)\right).$$

In this section we denote by $\nu = \nu_1 \otimes \nu_2$ the product distribution that generates the rewards from $\nu_j$ when pulling arm $j \in \{1, 2\}$. The regret of a policy that observes such rewards is denoted by $R_n(\nu)$. Finally let $\mathbb{P}_\nu$ denote the probability associated to $\nu$ and by $\mathbb{E}_\nu$ the corresponding expectation.

Hereafter, we favor rewards that are normally distributed because they lead to simpler calculations of the KL-divergence. However, our lower bounds remain of the same order for all families of distributions $\{\rho_\mu\}_\mu$ with expected value $\mu$ and such that $\mathrm{KL}(\rho_\mu - \rho_{\mu'}) \geq C(\mu - \mu')^2$ for some absolute constant $C > 0$. This is the case, for example, of the Bernoulli distribution with parameter $\mu$ as long as $\mu$ remains bounded away from 0 and 1; see, e.g., [11, Lemma 4.1].

The first lower bound illustrates that when one knows the distributions up to a permutation, the best one can hope for is a bounded regret of order $1/\Delta$.

THEOREM 5. *Let $\nu = \mathcal{N}(0, 1) \otimes \mathcal{N}(-\Delta, 1)$ and $\nu' = \mathcal{N}(-\Delta, 1) \otimes \mathcal{N}(0, 1)$. Then for any policy, and for every $n \geq 1$,*

$$\max\left(R_n(\nu), R_n(\nu')\right) \geq \frac{1}{4\Delta}.$$

PROOF. In this proof we assume that the policy has access to $t$ rewards from each arm at time step $t$. Clearly this full information setting is simpler than the bandit setting, and thus a lower bound for the former implies one for the latter. Using Lemma 4 as well as straightforward computations one obtains

$$\max\left(R_n(\nu), R_n(\nu')\right) \geq \frac{1}{2}\left(R_n(\nu) + R_n(\nu')\right) = \frac{\Delta}{2}\sum_{t=1}^{n}\left(\mathbb{P}_\nu(I_t = 2) + \mathbb{P}_{\nu'}(I_t = 1)\right)$$

$$\geq \frac{\Delta}{4}\sum_{t=1}^{n}\exp(-\mathrm{KL}(\nu^{\otimes t}, \nu'^{\otimes t})) = \frac{\Delta}{4}\sum_{t=1}^{n}\exp(-t\Delta^2) \geq \frac{1}{4\Delta}.$$

$\square$

The above theorem ensures that the regret bound of Theorem 1 has the correct dependence in $\Delta$. This is quite surprising as the original bound of [9] indicates that without the knowledge of $\mu^{(\star)}$ and $\Delta$, one can incur a regret that diverges to infinity at a logarithmic rate. The next result shows that this logarithmic regret already appears when one does not know the value of $\mu^{(\star)}$. Thus the knowledge of $\Delta$ without the knowledge of $\mu^{(\star)}$ is not sufficient to obtain a bounded regret. Moreover, the following lower bound matches the upper bounds (for the two-armed case) of [2], [4] and [10], thus proving their optimality.

THEOREM 6. *Let $\nu = \delta_0 \otimes \mathcal{N}(-\Delta, 1)$ and $\nu' = \delta_0 \otimes \mathcal{N}(\Delta, 1)$. Then for any policy, and any $n \geq 1$,*

$$\max\left(R_n(\nu), R_n(\nu')\right) \geq \frac{\log(n\Delta^2/2)}{4\Delta}.$$

PROOF. First note that

$$\max\left(R_n(\nu), R_n(\nu')\right) \geq R_n(\nu) \geq \Delta \mathbb{E}_\nu T_2(n).$$

Furthermore, denoting by $\nu_t$ (respectively $\nu'_t$) the law of the observed rewards up to time $t$ under $\nu$ (respectively under $\nu'$), and following the same computations than in the previous proof, one also obtains

$$\max\left(R_n(\nu), R_n(\nu')\right) \geq \frac{\Delta}{4} \sum_{t=1}^{n} \exp(-\mathrm{KL}(\nu_t, \nu'_t)).$$

Since under $\nu$, arm 1 is uninformative, it follows from basic calculation that

$$\mathrm{KL}(\nu_t, \nu'_t) = 2\Delta^2 \mathbb{E}_\nu T_2(t).$$

The above three displays yield

$$\begin{aligned}
\max\left(R_n(\nu), R_n(\nu')\right) &\geq \frac{\Delta}{2}\left(\mathbb{E}_\nu T_2(n) + \frac{n}{4}\exp(-2\Delta^2 \mathbb{E}_\nu T_2(n))\right) \\
&\geq \min_{x \in [0,n]} \frac{\Delta}{2}\left(x + \frac{n}{4}\exp(-2\Delta^2 x)\right) \\
&\geq \frac{\log(n\Delta^2/2)}{4\Delta}.
\end{aligned}$$

$\square$

Finally we prove that the knowledge of $\mu^{(\star)}$ without the knowledge of $\Delta$ is not sufficient either to obtain a bounded rescaled regret $\Delta R_n$. This result is more difficult, and falls within the more general topic of lower bounds for adaptive rates. First we need to generalize Lemma 4 to deal with both a composite alternative, and a rescaled risk. The proof of this result is standard and postponed to the appendix.

LEMMA 7. *Let $\rho_0$ and $\rho_\Delta, \Delta \in \mathbb{R}$ be probability distributions supported on some set $\mathcal{X}$, with $\rho_\Delta$ absolutely continuous with respect to $\rho_0$. Let $\lambda$ be a finite positive measure on $\mathbb{R}$. Then for any measurable function $\psi : \mathcal{X} \to \{0,1\}$, one has*

$$\mathbb{P}_{X \sim \rho_0}(\psi(X) = 1) + \int \Delta \mathbb{P}_{X \sim \rho_\Delta}(\psi(X) = 0) d\lambda(\Delta) \geq \frac{1}{C_\lambda} \exp\left(-\mathrm{KL}\left(\rho_0, \bar{\rho}\right)\right),$$

*where $\bar{\rho}$ is the positive measure on $\mathcal{X}$ defined by $\bar{\rho} = \int \Delta \rho_\Delta d\lambda(\Delta)$ and $C_\lambda = 1 + \int \Delta d\lambda(\Delta)$.*

Note that $\int \Delta \rho_\Delta \lambda(\Delta)$ is not a probability distribution, however it is a positive measure thus the Kullback-Leibler divergence in the above lemma is well-defined.

THEOREM 8. *Let $\nu_0 = \mathcal{N}(0,1) \otimes \mathcal{N}(-1,1)$, and $\nu_\Delta = \mathcal{N}(-\Delta,1) \otimes \mathcal{N}(0,1)$, $\Delta \in (0,1]$. Then for any policy, and any $n \geq 1$,*

$$\max\left(R_n(\nu_0), \sup_{\Delta \in (0,1]} \Delta R_n(\nu_\Delta)\right) \geq \frac{1}{2} \log(n/139).$$

Theorem 8 can be read as follows: for any policy, and any $n \geq 1$, there exists $\Delta \in (0,1]$ and a problem instance with gap $\Delta$ and optimal value $\mu^{(\star)} = 0$ such that on this problem one has

$$R_n \geq \frac{\log(n/139)}{2\Delta}.$$

PROOF. Similarly to the previous proof we define $\nu_{0,t}$ and $\nu_{\Delta,t}$ as the law of the observed rewards up to time $t$. Lemma 7 yields
(4.8)
$$\max\left(R_n(\nu_0), \sup_{\Delta \in (0,1]} \Delta R_n(\nu_\Delta)\right) \geq \frac{1}{2C_\lambda} \sum_{t=1}^{n} \exp\left(-\mathrm{KL}\left(\nu_{0,t}, \int \Delta \nu_{\Delta,t} d\lambda(\Delta)\right)\right).$$

For $\nu \in \{\nu_0, \nu_\Delta\}$, define the average rewards for arm $i \in \{1,2\}$ by $\mu_\nu^{(i)}$. Therefore, $\mu_{\nu_0}^{(1)} = \mu_{\nu_\Delta}^{(2)} = 0$, $\mu_{\nu_0}^{(2)} = -1$ and $\mu_{\nu_\Delta}^{(1)} = -\Delta$. Recall that a policy $\{I_t\}_{t \geq 1}$ taking values in $\{1,2\}$ generates a sequence of rewards $Y_t^{(I_t)}, t \geq 1$ distributed according to $\nu \in \{\nu_0, \nu_\Delta\}$. The joint density (with respect to the Lebesgue measure) $d\nu_t$ of $(Y_1^{(I_t)}, \ldots, Y_t^{(I_t)}) \in \mathbb{R}^t$, where $\nu \in \{\nu_\Delta, \nu_0\}$ can be computed easily using the chain rule for conditional densities. It is given by

$$d\nu_t = \frac{1}{(2\pi)^{t/2}} \exp\left(-\frac{1}{2} \sum_{\ell=1}^{t} (Y_\ell^{(I_\ell)} - \mu_\nu^{(I_\ell)})^2\right).$$

Choosing $\nu = \nu_\Delta$ and $\nu = \nu_0$ respectively, it yields

$$\frac{d\nu_{\Delta,t}}{d\nu_{0,t}}(Y_1^{(I_1)}, \ldots, Y_t^{(I_t)}) = \exp\left(-\frac{1}{2} \sum_{\ell=1}^{t} \left[(Y_\ell^{(I_\ell)} - \mu_{\nu_\Delta}^{(I_\ell)})^2 - (Y_\ell^{(I_\ell)} - \mu_{\nu_0}^{(I_\ell)})^2\right]\right)$$

$$= \exp\left(-\frac{1}{2} \sum_{\substack{\ell=1 \\ I_\ell=1}}^{t} \left[(Y_\ell^{(1)} + \Delta)^2 - (Y_\ell^{(1)})^2\right] - \frac{1}{2} \sum_{\substack{\ell=1 \\ I_\ell=2}}^{t} \left[(Y_\ell^{(2)})^2 - (Y_\ell^{(2)} + 1)^2\right]\right)$$

$$= \exp\left(-\frac{T^{(1)}}{2}(2\Delta\hat{\mu}^{(1)} + \Delta^2) + \frac{T^{(2)}}{2}(2\hat{\mu}^{(2)} + 1)\right),$$

where we denote for simplicity

$$T^{(i)} = T_i(t+1) = \sum_{\ell=1}^t \mathbb{1}\{I_\ell = i\} \quad \text{and} \quad \hat{\mu}^{(i)} = \hat{\mu}^{(i)}_{T_i(t)} = \frac{1}{T^{(i)}} \sum_{\substack{\ell=1 \\ I_\ell=i}}^t Y_\ell^{(i)}, \quad i \in \{1,2\}.$$

Dropping the dependency in $(Y_1^{(I_1)}, \ldots, Y_t^{(I_t)})$ from the notation, it yields

$$\int \Delta \frac{d\nu_{\Delta,t}}{d\nu_{0,t}} d\lambda(\Delta) = \exp\left(\frac{T^{(2)}}{2}(2\hat{\mu}^{(2)} + 1)\right) \int \Delta \exp\left(-\frac{T^{(1)}}{2}(2\Delta\hat{\mu}^{(1)} + \Delta^2)\right) d\lambda(\Delta),$$

and thus

$$\text{KL}\left(\nu_{0,t}, \int \Delta \nu_{\Delta,t} d\lambda(\Delta)\right)$$

$$= -\mathbb{E}_{\nu_0}\left(\frac{T^{(2)}}{2}(2\hat{\mu}^{(2)} + 1) + \log\left(\int \Delta \exp\left(-\frac{T^{(1)}}{2}(2\Delta\hat{\mu}^{(1)} + \Delta^2)\right) d\lambda(\Delta)\right)\right)$$

$$= \frac{1}{2}\mathbb{E}_{\nu_0} T^{(2)} - \mathbb{E}_{\nu_0} \log\left(\int \Delta \exp\left(-\frac{T^{(1)}}{2}(2\Delta\hat{\mu}^{(1)} + \Delta^2)\right) d\lambda(\Delta)\right)$$

where the last line follows standard computations. Next, it follows from the Cauchy-Schwarz inequality that the function

$$x \mapsto \log\left(\int_\Delta \Delta \exp(\varphi(\Delta)x) d\lambda(\Delta)\right)$$

is convex for any function $\varphi$. Together with the Jensen inequality, it yields

$$\mathbb{E}_{\nu_0} \log\left(\int \Delta \exp\left(-\frac{T^{(1)}}{2}(2\Delta\hat{\mu}^{(1)} + \Delta^2)\right) d\lambda(\Delta)\right)$$

$$\geq \log\left(\int \Delta \exp\left(-\mathbb{E}_{\nu_0}\frac{T^{(1)}}{2}(2\Delta\hat{\mu}^{(1)} + \Delta^2)\right) d\lambda(\Delta)\right)$$

$$= \log\left(\int \Delta \exp\left(-\frac{\mathbb{E}_{\nu_0} T^{(1)}}{2}\Delta^2\right) d\lambda(\Delta)\right)$$

Define $\tau = \mathbb{E}_{\nu_0} T^{(1)}$ and let $\lambda$ be the uniform distribution on $[0, 1/\sqrt{\tau}]$. Since $ue^{-u^2/2} \geq u/2$ for $0 \leq u \leq 1$, it yields

$$\int \Delta \exp\left(-\frac{\mathbb{E}_{\nu_0} T^{(1)}}{2}\Delta^2\right) d\lambda(\Delta) = \frac{1}{\sqrt{\tau}}\int_0^1 u \exp(-u^2/2) du \geq \frac{1}{4\sqrt{\tau}},$$

Thus we have proved that

$$\text{KL}\left(\nu_{0,t}, \int_0^1 \Delta \nu_{\Delta,t} d\Delta\right) \leq \frac{1}{2}\mathbb{E}_{\nu_0} T^{(2)} + \log(4\sqrt{\mathbb{E}_{\nu_0} T^{(1)}})$$

$$\leq \frac{1}{2}\mathbb{E}_{\nu_0} T_2(n) + \frac{1}{2}\log(16n).$$

Plugging this into (4.8) one obtains

$$\max\left(R_n(\nu_0), \sup_{\Delta \in (0,1]} \Delta R_n(\nu_\Delta)\right) \geq \frac{\sqrt{n}}{8C_\lambda} \exp\left(-\frac{1}{2}\mathbb{E}_{\nu_0}T_2(n)\right)$$

$$\geq \frac{\sqrt{n}}{16} \exp\left(-\frac{1}{2}\mathbb{E}_{\nu_0}T_2(n)\right),$$

where we use the fact that $\tau \geq 1$, which implies $C_\lambda \leq 3/2 \leq 2$. On the other hand one also has

$$R_n(\nu_0) \geq \mathbb{E}_{\nu_0}T_2(n)$$

Therefore

$$\max\left(R_n(\nu_0), \sup_{\Delta \in (0,1]} \Delta R_n(\nu_\Delta)\right) \geq \min_{x \in [0,n]} \frac{1}{2}\left(x + \frac{\sqrt{n}}{16}\exp(-x/2)\right)$$

$$= \frac{1}{2}\log(n/139).$$

$\square$

Theorem 6 and 8 have important consequences on the *exploration-exploitation tradeoff* mentioned in the introduction. Indeed, consider the full information case where at each round, the agent observes the reward of both arms. In this case, it is not hard to see that the policy that indicates to pull the arm with the best average reward has bounded regret of order $1/\Delta$. Therefore, the knowledge of $\Delta$ or $\mu^{(\star)}$ alone does not alleviate the price for exploration. However, when both are known, it vanishes (see Theorem 1).

## APPENDIX A: PROOF OF LEMMA 7

Throughout the proof, Radon-Nikodym derivatives over $\mathcal{X}$ are taken with respect to a common but unspecified reference measure. It does not enter our final result. It follows from Fubini's Theorem that

$$\mathbb{P}_{X \sim \rho_0}(\psi(X) = 1) + \int \Delta \mathbb{P}_{X \sim \rho_\Delta}(\psi(X) = 0)d\lambda(\Delta)$$

$$= \int_{\psi=1} d\rho_0 + \int \left(\int_{\psi=0} \Delta d\rho_\Delta\right) d\lambda(\Delta)$$

$$= \int_{\psi=1} d\rho_0 + \int_{\psi=1} d\bar{\rho}$$

$$= \int_{\psi=0} d\rho_0 + \int_{\psi=1} \frac{d\bar{\rho}}{d\rho_0} d\rho_0$$

Furthermore the last expression is clearly minimized for $\psi(x) = \mathbb{1}\left\{\frac{d\bar{\rho}}{d\rho_0}(x) > 1\right\}$.
It yields

$$
\int_{\psi=1} d\rho_0 + \int_{\psi=0} \frac{d\bar{\rho}}{d\rho_0} d\rho_0 \geq \int_{\frac{d\bar{\rho}}{d\rho_0}>1} d\rho_0 + \int_{\frac{d\bar{\rho}}{d\rho_0}\leq 1} \frac{d\bar{\rho}}{d\rho_0} d\rho_0(x)
$$

$$
= \int_{\frac{d\bar{\rho}}{d\rho_0}>1} d\rho_0 + \int_{\frac{d\bar{\rho}}{d\rho_0}\leq 1} d\bar{\rho}
$$

$$
= \int \min\left(d\rho_0, d\bar{\rho}\right).
$$

Note that the latter quantity is often referred to as *Hellinger affinity* and does not depend on the reference measure on $\mathcal{X}$; see, e.g., [14], Chapter 2. Now using the Cauchy-Schwarz inequality and the fact that

$$
\int \min\left(d\rho_0, d\bar{\rho}\right) + \int \max\left(d\rho_0, d\bar{\rho}\right) = C_\lambda,
$$

we get

$$
\left(\int \sqrt{d\bar{\rho} d\rho_0}\right)^2 = \left(\int \sqrt{\min(d\bar{\rho}, d\rho_0)\max(d\bar{\rho}, d\rho_0)}\right)^2
$$

$$
\leq \left(\int_x \min(d\bar{\rho}, d\rho_0)\right)\left(\int_x \max(d\bar{\rho}, d\rho_0)\right)
$$

$$
\leq C_\lambda \int_x \min(d\bar{\rho}, d\rho_0).
$$

The above three displays together yield

$$
\mathbb{P}_{X\sim\rho_0}(\psi(X) = 1) + \int_\Delta \Delta \mathbb{P}_{X\sim\rho_\Delta}(\psi(X) = 0)d\lambda(\Delta) \geq \frac{1}{C_\lambda}\left(\int \sqrt{d\bar{\rho} d\rho_0}\right)^2.
$$

To complete the proof, observe that the Jensen inequality yields

$$
\left(\int \sqrt{d\bar{\rho} d\rho_0}\right)^2 = \left(\int \sqrt{\frac{d\bar{\rho}}{d\rho_0}} d\rho_0\right)^2
$$

$$
= \exp\left[2\log\left(\int \sqrt{\frac{d\bar{\rho}}{d\rho_0}} d\rho_0\right)\right]
$$

$$
\geq \exp\left[2\int \log\left(\sqrt{\frac{d\bar{\rho}}{d\rho_0}}\right) d\rho_0\right]
$$

$$
= \exp[-\mathrm{KL}(\rho_0, \bar{\rho})].
$$

## REFERENCES

[1] Agrawal, R., Teneketzis, D., and Anantharam, V. Asymptotically efficient adaptive allocation schemes for controlled i.i.d. processes: finite parameter space. *IEEE Trans. Automat. Control 34*, 3 (1989), 258–267.

[2] Audibert, J.-Y., and Bubeck, S. Minimax policies for adversarial and stochastic bandits. In *Proceedings of the 22nd Annual Conference on Learning Theory (COLT)* (2009).

[3] Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine Learning Journal 47*, 2-3 (2002), 235–256.

[4] Auer, P., and Ortner, R. UCB revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica 61*, 1 (2010), 55–65.

[5] Bubeck, S., and Cesa-Bianchi, N. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning 5*, 1 (2012), 1–122.

[6] Kulkarni, S. R., and Lugosi, G. Finite-time lower bounds for the two-armed bandit problem. *IEEE Transactions on Automatic Control 45*, 4 (2000), 711–714.

[7] Lai, T. L., and Robbins, H. Asymptotically optimal allocation of treatments in sequential experiments. In *Design of Experiments: Ranking and Selection*, T. J. Santner and A. C. Tamhane, Eds. 1984, pp. 127–142.

[8] Lai, T. L., and Robbins, H. Optimal sequential sampling from two populations. *Proc. Natl. Acad. Sci. USA 81* (1984), 1284–1286.

[9] Lai, T. L., and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics 6* (1985), 4–22.

[10] Perchet, V., and Rigollet, P. The multi-armed bandit problem with covariates, October 2011. arXiv:1110.6084.

[11] Rigollet, P., and Zeevi, A. Nonparametric bandits with covariates. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)* (2010), A. T. Kalai and M. Mohri, Eds., pp. 54–66.

[12] Robbins, H. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematics Society 58* (1952), 527–535.

[13] Salomon, A., and Audibert, J.-Y. Deviations of stochastic bandit regret. In *Proceedings of the 22nd International Conference on Algorithmic Learning Theory (ALT)* (2011).

[14] Tsyabkov, A. B. *Introduction to Nonparametric Estimation*. Springer, 2009.

Sébastien Bubeck
Department of Operations Research
 and Financial Engineering
Princeton University
Princeton, NJ 08544, USA
(sbubeck@princeton.edu)

Vianney Perchet
LPMA, UMR 7599
Université Paris Diderot
175, rue du Chevaleret
75013 Paris, France
(vianney.perchet@normalesup.org)

Philippe Rigollet
Department of Operations Research
 and Financial Engineering
Princeton University
Princeton, NJ 08544, USA
(rigollet@princeton.edu)